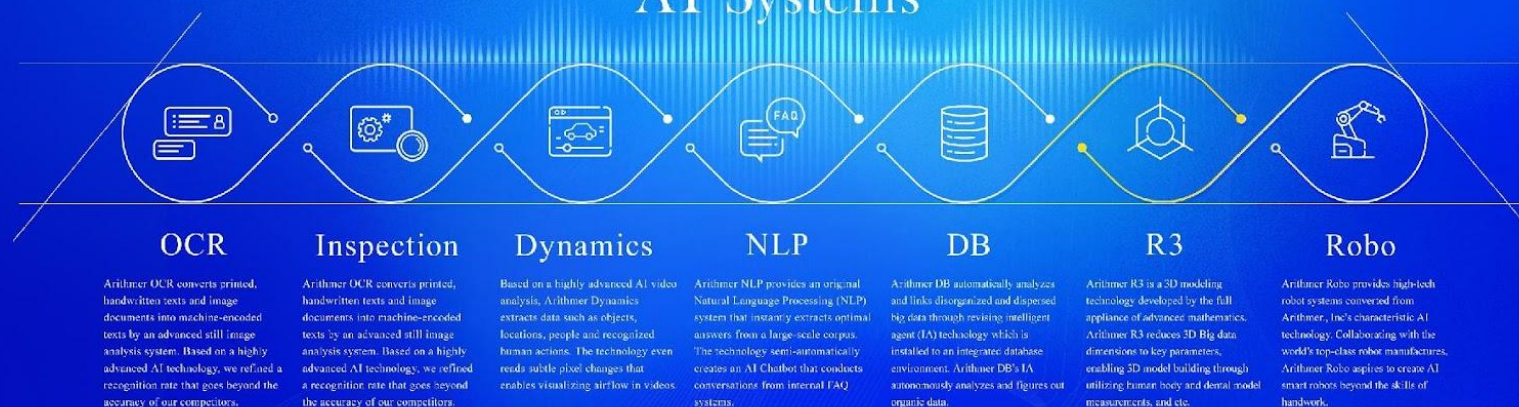


Arithmer R3

AI Systems



Weakly Supervised Semantic Segmentation in 3D Graph-Structured Point Clouds of Wild Scenes

Arithmer R3 team
Daisuke SATO

2021/7/18

<https://arxiv.org/abs/2004.12498v2>

Point clouds

- 3D sensors have been developed rapidly these days.
 - iPhone/iPad has LiDAR
 - realsense/azure kinect is amazing considering its low price.
- Raw data collected is point cloud
 - Collection of 3D points
 - (x_0, y_0, z_0)
 - (x_1, y_1, z_1)
 - ...
 - (x_n, y_n, z_n)
- Processing point cloud is quite important in robotics/measurement (測量).



3D Point Cloud



Task: semantic segmentation of point clouds

- Classifying every point of 3D point clouds
 - If we perform this task with supervised learning, we need to annotate each point.
 - **Labeling point clouds is super exhausting.**
- Labeling 2D image is easier.
- **Let us do weakly supervised learning with 2D projected image.**



3D Point Cloud



3D Segmentation Output



Idea: weakly supervised learning with 2D projected image

Input is
3D point cloud

3D Point Cloud



Viewpoints

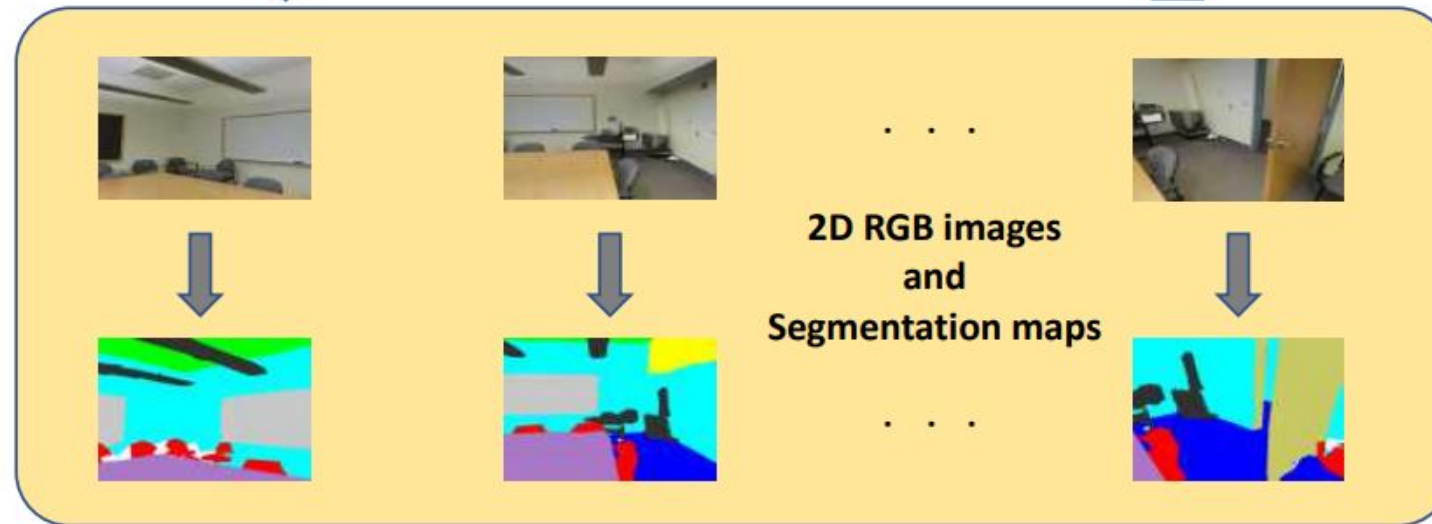
3D Segmentation Output



Supervision



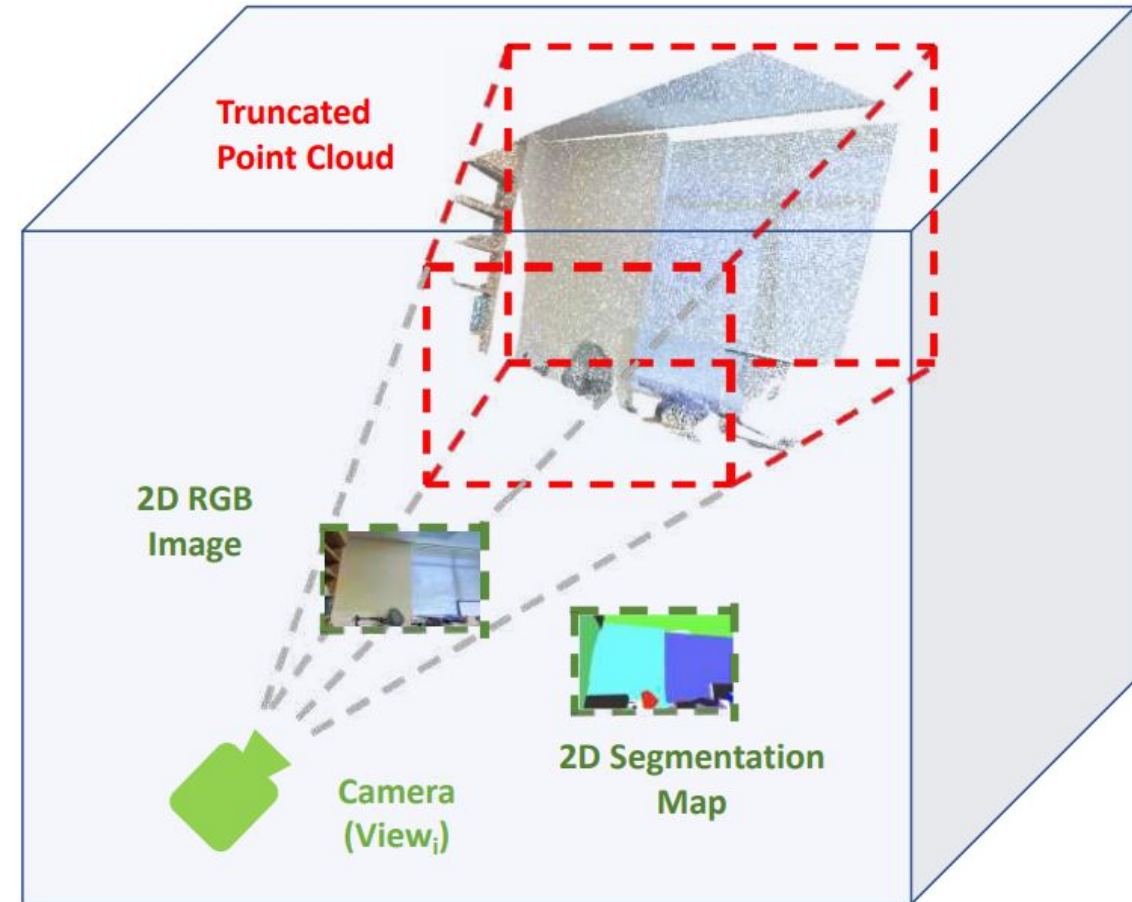
Output at
inference stage
is 3D segmented
point cloud



Label at training
stage is
segmented 2D
projected images

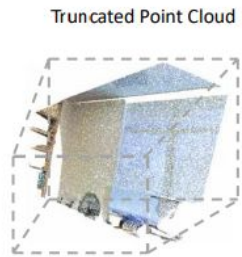
Note: Projection of 3D pointcloud onto 2D image

- Projection is described by the two parameters:
 - Internal camera parameters
 - Focal length (f)
 - principal point (c)
 - External camera parameters
 - Translation (t)
 - Rotation (R)



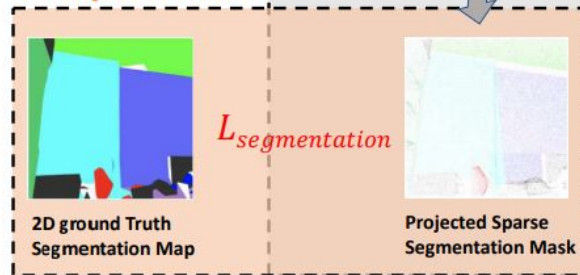
Model architecture

Input is 3D point cloud



2D projection
(segmented pcd/mask for visibility
-> 2D segmented image)

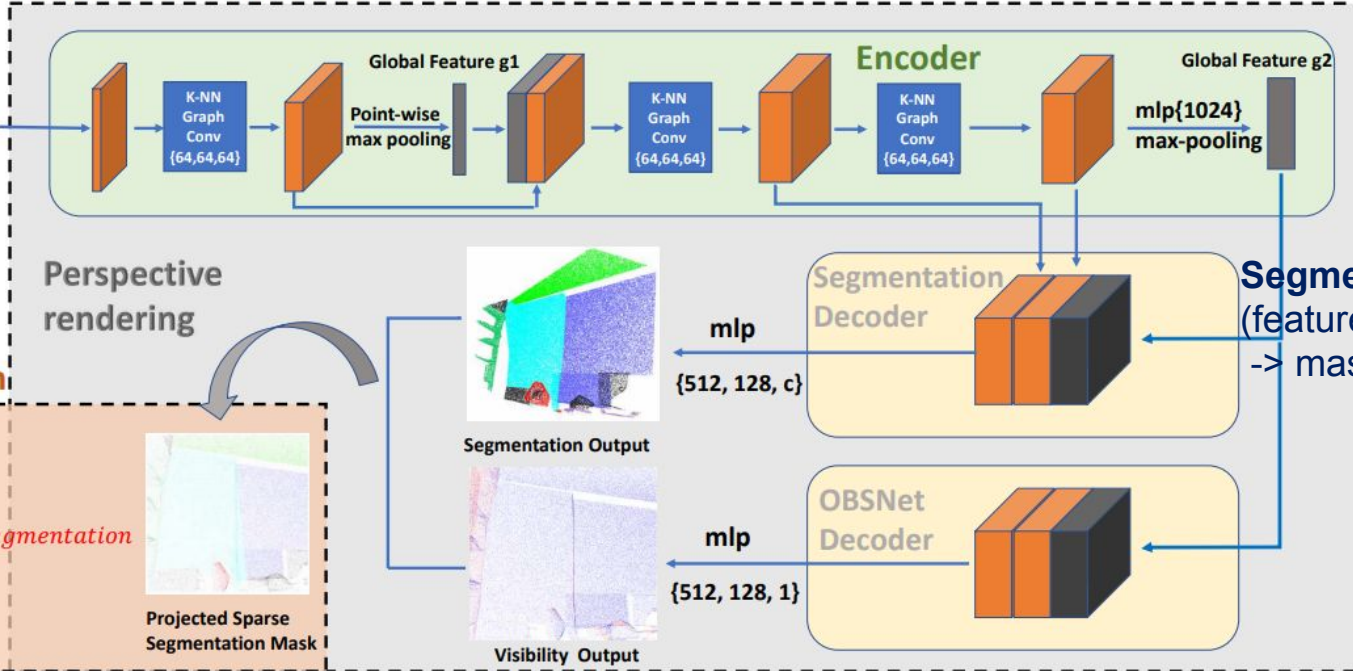
2D Optimization



$$L_{seg} = -\frac{1}{N} \sum_{i=1}^N [p_i \log \hat{p}_i + (1 - p_i) \log(1 - \hat{p}_i)]$$

Total loss $L = L_{seg} + \lambda L_{vis}$

Graph-based Pyramid Feature Network (GPFN)



Graph convolution encoder
(point cloud -> feature vector)

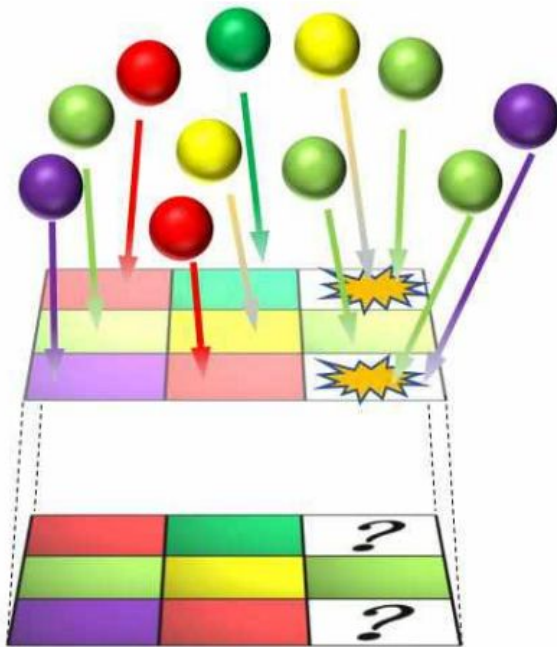
Segmentation branch
(feature vector
-> mask for segmentation)

Visibility branch
(feature vector -> mask for visibility)

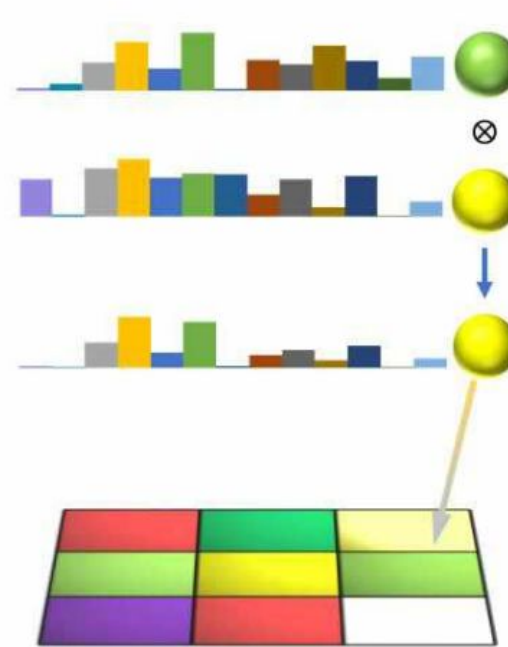
$$L_{vis} = -\frac{1}{M} \sum_{i=1}^M [U_i \log \hat{U}_i + (1 - U_i) \log(1 - \hat{U}_i)]$$

Perspective rendering

- Problem: Multiple pointclouds can be projected onto a single pixel. Which class should we give the pixel?



- Solution: “Semantic fusion” (sophiscated voting system)



$$p(C_i|x_{grid}) = \prod_{n=1}^N p(C_i|x_n),$$

$$p(C_i|x_{grid})_{norm} = p(C_i|x_{grid}) / \sum_{i=1}^{n_{classes}} \prod_{n=1}^N p(C_i|x_n),$$

$$p(x_{grid}) = \max\{p(C_1|x_{grid}), \dots, p(C_{n_{classes}}|x_{grid})\}.$$

Datasets

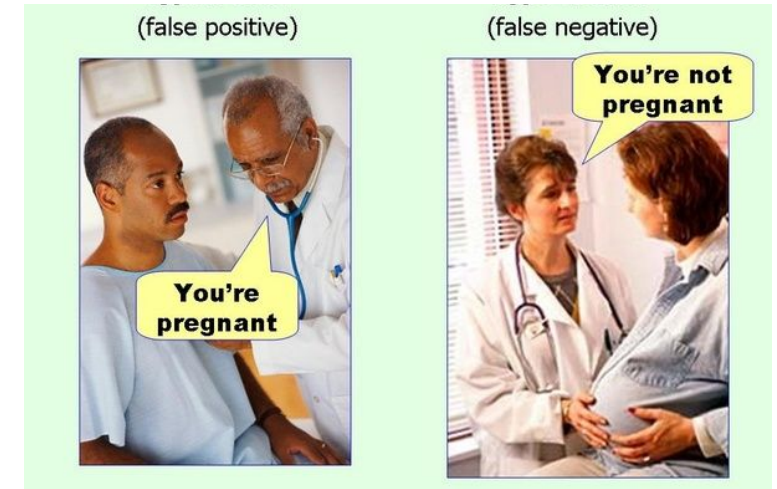
- **SUNCG Synthetic** dataset
 - class number: 40 (furniture)
 - rooms: 404058
 - create 55000 2D rendering sets
- **S3DIS Real-world** dataset
 - class number: 13 (furniture)
 - rooms: 272
 - thousands of viewpoints are provided
- Each point has
 - position: (x, y, z)
 - color: (r, g, b)
- Also normal is computed (u, v, w)



Metrics

- mean accuracy of total classes ($mAcc$)
- overall accuracy ($oAcc$)
- mean per-class intersection-over-union ($mIoU$)
 - $IOU = TP / (TP + FP + FN)$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$



$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

The diagram shows two Venn diagrams. The top one has two white circles with a grey shaded intersection, representing the 'Area of Overlap'. The bottom one has two grey circles with a grey shaded intersection, representing the 'Area of Union'.

Comparison with other fully-supervised methods

	Method	mAcc(%)	mIoU(%)	oAcc(%)
3D Supervision	<i>PointNet</i> [28]	66.2	47.6	78.5
	<i>Engelmann et al.</i> [10]	66.4	49.7	81.1
	<i>PointNet++</i> [29]	67.1	54.5	81.0
	<i>DGCNN</i> [40]	-	56.1	84.1
	<i>Engelmann et al.</i> [11]	67.8	58.3	84.0
	<i>SPG</i> [19]	73.0	62.1	85.5
2D Supervision	<i>GPFN with DP (Ours)</i>	39.2	30.4	53.7
	<i>GPFN with DP w/ D_v (Ours)</i>	59.4	42.7	70.0
	<i>GPFN with PR w/o D_v (Ours)</i>	54.2	39.0	66.8
	<i>GPFN with PR w/ D_v (Ours)</i>	66.5	50.8	79.1

Not so bad even compared with fully-supervised learning.

Inference samples for SUNCG dataset

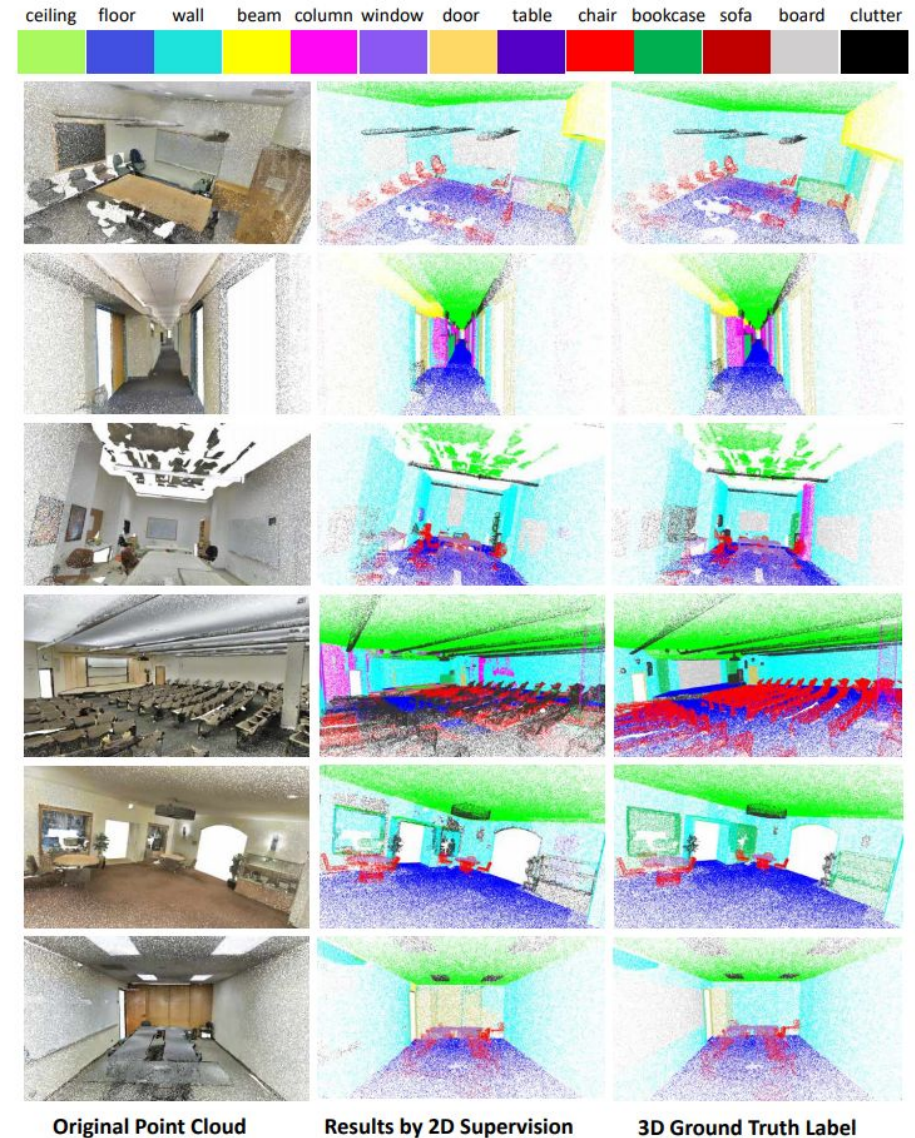


Original Point Cloud

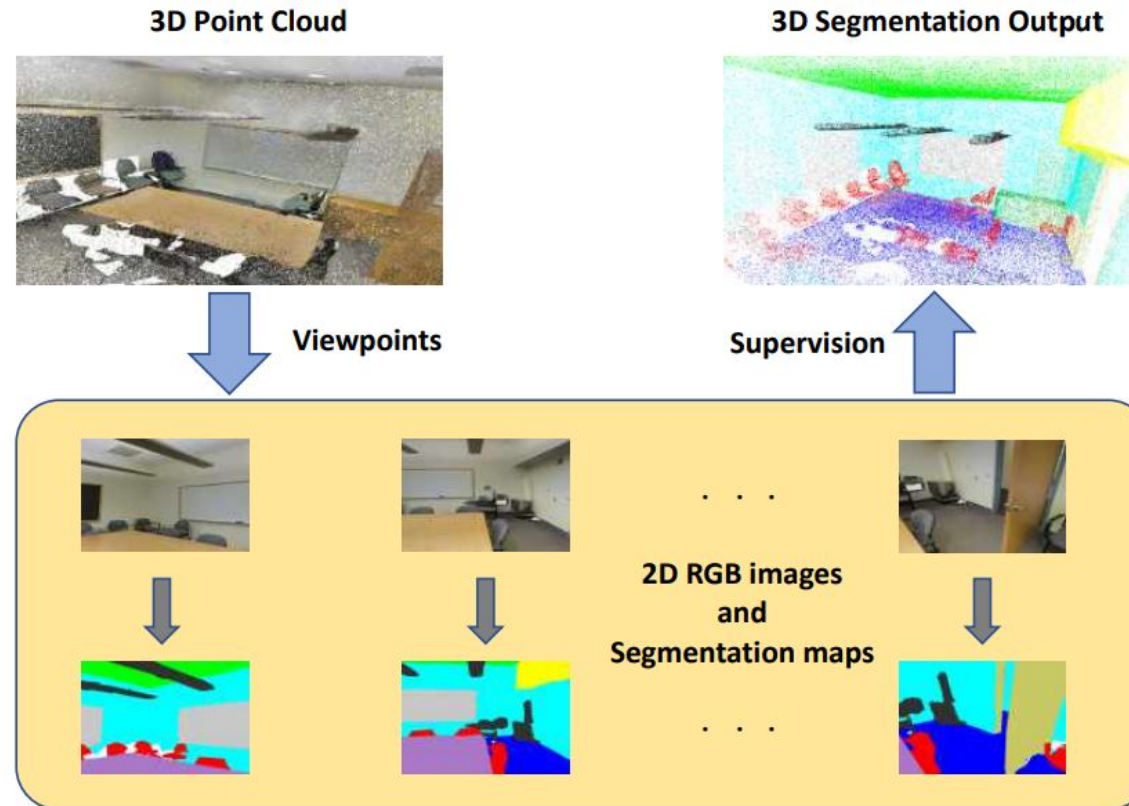
Results by 2D Supervision

3D Ground Truth Label

Inference samples for S3DIS dataset



- To the best of our knowledge, this is the **first work to apply 2D supervision for 3D semantic point cloud segmentation of wild scenes without using any 3D pointwise annotations.**
- Extensive experiments are conducted and the proposed method achieves comparable performance with the state-of-the-art 3D supervised methods on the popular SUNCG and S3DIS benchmarks.



人間に、愛を。
未来に、AIを。

Arithmer 株式会社

〒106-6040

東京都港区六本木一丁目6番1号 泉ガーデンタワー 38/40F(受付)

03-5579-6683

<https://arithmer.co.jp/>

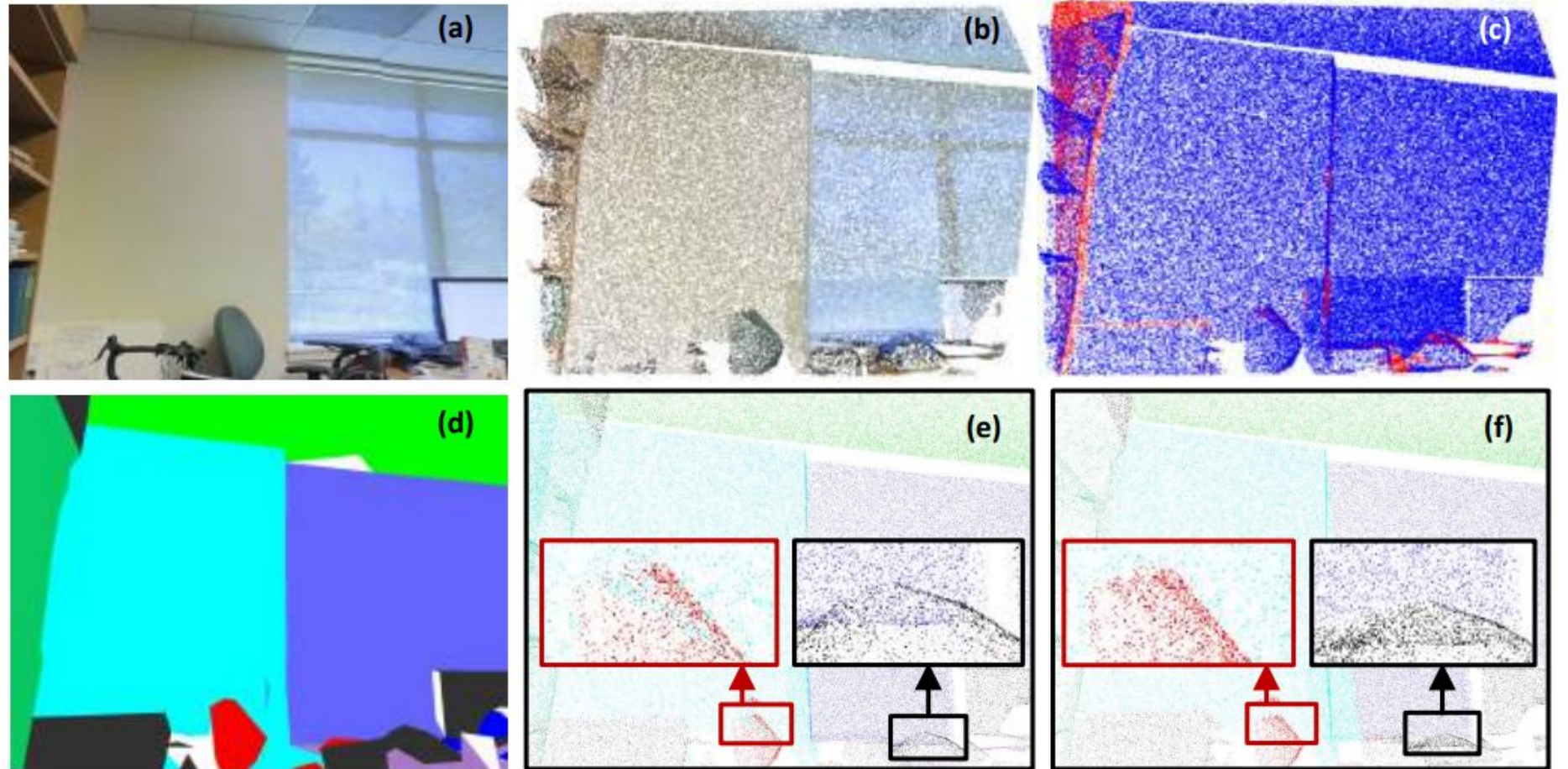
Arithmer



BACK UP

TITLE HERE

- Bullet list



Ablation study: projection method/OBSNet decoder

- Bullet list

Method	mAcc(%)	mIoU(%)	oAcc(%)
<i>GPFN with DP (Ours)</i>	61.9	45.0	73.4
<i>GPFN with DP w/ D_v (Ours)</i>	71.9	61.2	84.5
<i>GPFN with PR w/o D_v (Ours)</i>	65.3	50.8	79.1
<i>GPFN with PR w/ D_v (Ours)</i>	87.3	70.37	91.8

Ablation study - Encoder Design

- Bullet list

K-NN Graph	Pyramid	mAcc(%)	mIoU(%)	oAcc(%)
×	×	61.3	45.1	72.6
✓	×	65.1	48.6	78.4
×	✓	63.5	46.4	75.3
✓	✓	66.5	50.8	79.1

Ablation study - Amount of training data

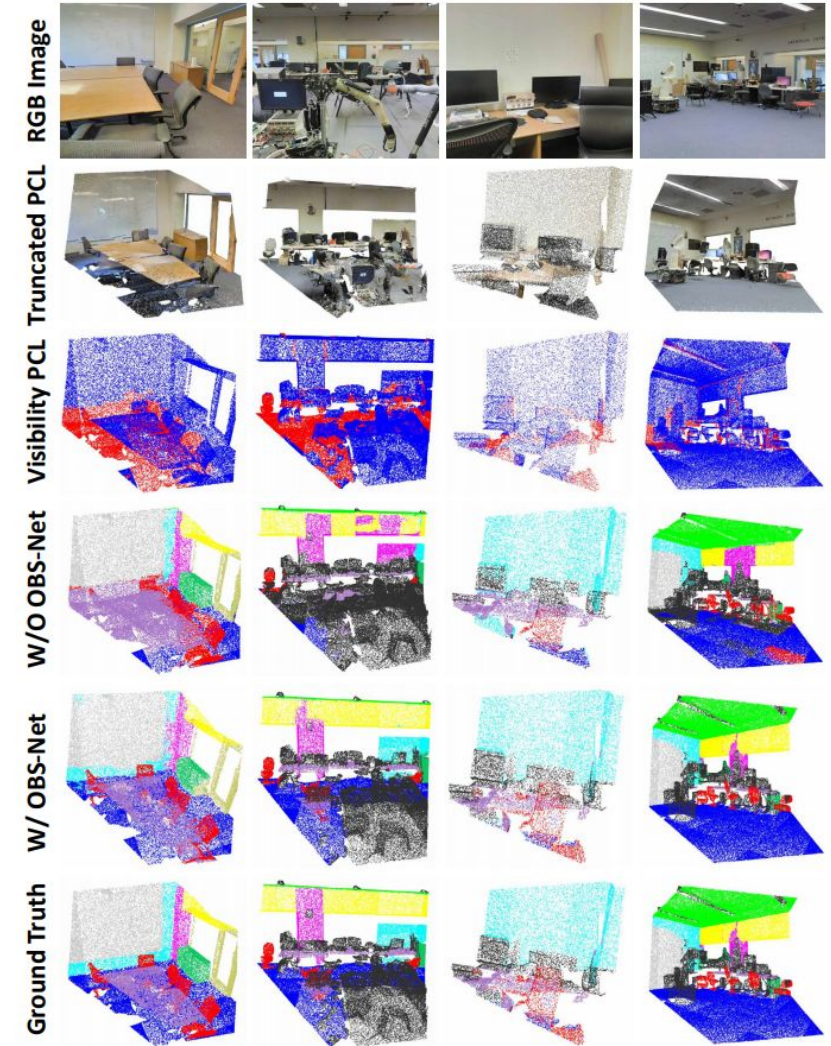
- Bullet list

Training data	mAcc (%)	mIoU (%)	oAcc (%)
<i>All</i>	67.0	52.5	81.5
1/2	66.9	51.8	80.9
1/4	66.7	50.9	79.5
1/6	66.5	50.8	79.1
1/12	56.5	39.3	66.2
1/20	37.8	29.1	40.0

Ablation study - Visibility detection by OBSNet

- Bullet list

Dataset	Accuracy (%)					
	All	1/2	1/4	1/6	1/12	1/20
S3DIS	93.0	92.6	91.7	91.2	89.6	85.0



Transfer learning from synthetic to realistic dataset

- Bullet list

Training Data	mAcc(%)	mIoU(%)	oAcc(%)
Train Scratch on S3DIS	66.5%	50.8%	79.1%
Pretrained on SUNCG	67.0%	53.5%	81.3%